

Abstract Title Page
Not included in page count.

Title: What Drives Alignment between Teachers' Survey Self-Reports and Classroom Observations of Standards-Based Mathematics Instruction?

Author(s): Julia Kaufman & Brian Junker

Abstract Body

Background / Context: Research on whether teachers can give accurate self-reports in surveys about their mathematics instruction is fairly mixed. Some of that research indicates that teachers can provide some general approximation of their mathematics instruction in survey self-reports (Mayer 1999; Ross, McDougall et al. 2003), while other studies find very little correlation between teachers' survey responses and their mathematics practices (Spillane and Zeuli 1999; Stecher, Le et al. 2006). This diversity of findings points to elements within school district and program implementation context that influence teachers' understanding of mathematics instruction and, thus, the accuracy of their reports about that instruction. While some research provides evidence that teachers' understanding drives the accuracy of their survey responses (Hill 2005; Spillane and Zeuli 1999), no research to date has provided evidence about what aspects of district context impact the accuracy of teachers' self-reports and whether that accuracy can change over time.

In this paper, we consider the accuracy of teachers' survey reports about their mathematics instruction over a two-year period in two urban school districts. Our work suggests that several elements of district context matter for the accuracy of teachers' self reports, including their mathematics learning opportunities and the presence of other big instructional initiatives within the district. These findings, drawn from in-depth quantitative and qualitative data gathered in two localized education settings, provide key hypotheses to guide future survey research and controlled studies on teachers' understanding of their mathematics instructional practices (e.g., as suggested by Shavelson & Towne, 2002, pp 105-108).

In "Scaling Up Mathematics" – the five-year, federally-funded project from which we draw our data – researchers investigated the implementation of standards-based elementary school mathematics curricula in two school districts in order to understand how teachers' knowledge and social interaction impacted the quality of their instruction.[†] "Standards-based" mathematics curricula encourage the development of students' conceptual understanding through explorations and problem solving; many of these curricula sprang up in response to National Council of Teachers of Mathematics' *Curriculum and Evaluation Standards for School Mathematics* (1989). For Scaling Up Mathematics, teachers' standards-based instruction was measured through teacher surveys and observations.

Other Scaling Up Mathematics work has demonstrated that teachers in one district had considerably higher-quality standards-based instruction compared to the other (Stein and Kaufman 2010), and that higher-quality instruction could be tied to the stronger curriculum-based professional development and the more in-depth social interaction documented in that district (Stein and Kaufman 2010; Coburn and Russell 2008). In the present study, we find that the same district with higher-quality instruction also had many more teachers who provided accurate self-reports of that instruction, while teachers in the other district often overestimated the quality of their instruction in surveys. Using additional survey and interview data, we present some individual and district factors that explain the accuracy of teachers' responses.

Purpose / Objective / Research Question / Focus of Study: Our research questions are: 1.) How closely do teachers' reports of their standards-based pedagogical practices align with their actual standards-based pedagogical practices, as measured through classroom observations? 2.) How does the accuracy of teachers' reports of their practice change across the two years of our study? 3.) How

[†] Scaling Up Mathematics was supported by a grant from the Interagency Educational Research Initiative. Stein and Coburn (2008) compare the curriculum implementation process in the two districts.

do district context and teacher-level factors relate to the accuracy of teachers' responses regarding their standards-based practice?

Setting: We draw our data from the Scaling Up Mathematics project, which collected data about teachers' instruction in two large urban school districts implementing different standards-based mathematics curricula: *Everyday Mathematics* in Region Z in New York City and *Investigations* in Greene (district names are pseudonyms).

Population / Participants / Subjects: Participating teachers are from four case study schools in each district (eight schools altogether). Case study schools were chosen to represent a range in the extent of teacher knowledge and social interaction within each district. In each school, the principal nominated a teacher from each grade to participate in the study. After teacher attrition and replacing some teachers, our final participants for this study are 47 teachers, 23 in Region Z and 24 in Greene.

Intervention / Program / Practice: The goal of this study is to compare the validity of teacher responses relative to a gold standard – classroom observation – in two districts that differ with regard to the standards-based curricula that they are implementing and the learning opportunities that are provided to teachers. To do this work, we use 1.) quantitative data from surveys and classroom observation coding and 2.) qualitative data from observation field notes and teacher interviews to provide detailed, in-depth information about what is happening in two localized education settings (Shavelson & Towne, 2002, pp 105-108), Greene and Region Z.

Research Design: We use a comparative embedded case study approach (Yin 1994) to investigate standards-based instructional practices in two districts and four schools in each district as measured through 1.) 431 coded lesson observations (206 in Region Z and 225 in Greene) among 47 teachers and 2.) teachers' annual self-reports about their standards-based pedagogical practices. To provide explanations for variation in the accuracy of teachers' self-reports, we link our findings on quantitative patterns in the accuracy of teachers' responses with intensive interview data from eight teachers in four of the schools, two in each district.

Data Collection and Analysis: We have classroom observation data from at least two and most often six lessons per teacher each year, as well as one survey self-report each year. See Table 1 for a listing of all teacher participants, including their years of participation in the study, their number of classroom observations each year, and the accuracy of their survey response. We did not include teachers in the study if we did not have at least one classroom observation for each of semester (spring and fall) of the year alongside their survey response.

[Please insert Table 1 here.]

We collected a wide range of data in the lesson observations and surveys, including specific information about two key “standards-based” practices described in the mathematics education literature: 1.) teachers' work to uncover student thinking (Lampert 1990; Shifter 2001) and 2.) the extent to which teachers helped students use mathematics (versus the teacher or text) to justify their thinking (Lampert 1990; Engle and Faux 2006). For the remainder of this abstract, we call those two teacher practices “uncovering thinking” and “justifying with math.”

Surveys. Surveys for the Scaling Up Mathematics project were administered in the spring of each year to all teachers in both districts, including the subset of teachers for whom we have observation data. In addition to questions meant to capture details about teachers' experience, education, mathematics classes, and curriculum use, the survey included many questions about teachers'

standards-based practices.[‡] From the survey, we chose more than 30 items that measured uncovering thinking and justifying with math. Based on factor analysis of the first year data and further analysis of individual items, we decided to use 23 of those items in a single composite measuring uncovering thinking and justifying with math ($\alpha=.83$). We used a single composite – rather than two separate composites to measure uncovering thinking and justifying with math – because many of those items conflate teachers’ uncovering student thinking and justifying with math work.

Classroom Observations and Interviews. Trained observers recorded detailed notes about each lesson and also responded to specific prompts about uncovering thinking and justifying with math. Those observers also interviewed teachers four times each year – generally before and after a lesson observation – about their lesson plans, people with whom they spoke about their lesson plans, and their general thoughts on mathematics instruction.

Lesson observation notes were later coded by experienced mathematics educators who assigned a rating of 0, 1 or 2 for teachers’ work to uncover thinking and a rating of 0, 1, or 2 for teachers’ work to help students justify with mathematics. Two mathematics educators coded 43 lessons together; those coders gave teachers the same “uncovering student thinking” rating for 74% of the lessons and the same “justifying with math” rating for 79% of the lessons. In order to compare our survey composite of teachers’ standards-based practices to observations, we use the average of the two observer ratings for uncovering thinking and justifying with math as a single observation measure of teachers’ standards-based practices in one lesson, and we calculated a yearlong observation score as an average of all those lesson observations across a year for each teacher. For details on all our measures of standards-based practices – the 23 survey items, the prompts for observers and rating scale for coders of observations – see Table 2.

[Please insert Table 2 here.]

Finally, we used conceptually-determined cut scores for both the survey and observation measures of standards-based practice to define “high,” “medium, and “low” standards-based practice. For a “low” score, a teacher’s annual survey measure or observation rating reflects doing no or almost no uncovering thinking or justifying with math work; a “medium” score reflects some but not regular and frequent uncovering thinking and justifying with math work; and a “high” score reflects a teacher’s frequent, regular work to uncover thinking and justify with math.

Findings / Results:

Descriptive Data on Teachers’ Standards-Based Practice in Survey Self-Reports and Observations

Figures 1 and 2 below provide averages and confidence intervals for, respectively, teachers’ survey self-reports and their observation ratings of their standards-based pedagogical practices in the two school districts. As can be seen in the figures, survey reports are significantly higher for Greene teachers – compared to Region Z teachers – for Year 1 only. However, observation ratings are clearly higher for Greene teachers for both Year 1 and Year 2. These same trends – including significant differences between districts for the observations – are present for the smaller number of teachers for whom we have two years of data.

[Please insert Figures 1 and 2 about here.]

Relationship between Teachers’ Survey Self-Reports and Observation Ratings

We found positive and significant ($p<.01$) correlations between Greene teachers’ survey self-reports and observation ratings (.66 in Year 1 and .65 in Year 2). In contrast, correlations for Region Z teachers were much lower and not significant (.39 in Year 1 and .06 in Year 2). Figures 3 and 4 are

[‡] Some of our items are drawn from a survey measure by Ross, McDougall et al (2003) meant to capture the extent of teachers’ standards-based mathematics practices.

scatter plots illustrating the relationship between survey self-reports and observation ratings separately for Year 1 and Year 2. The blue lines in the plots represent cut scores for high, medium, and low standards-based instruction in the surveys and observations. These plots demonstrate that most Greene teachers provide accurate ratings of their practice – most with “medium” survey ratings and “medium” observation ratings – while many Region Z teachers overestimate their practices, with survey reports of “medium” standards-based pedagogical practices and “low” observation ratings. Generally, those teachers with “medium” standards-based practices accurately estimated the extent of those practices in a survey, whereas teachers with “high” or “low” practices were not as accurate.

[Please insert Figures 3 and 4 here.]

When teachers change from one year to the next, either their survey self-report changes by one step (e.g. from “low” to “medium”) or their observation rating changes by one step. These findings suggest that any changes to teachers’ survey self-report or practice are more likely to be incremental, rather than big leaps. Furthermore, change in survey self-reports may follow or precede change to actual practice, but change in self-reports and practice does not occur simultaneously.

Explanations for the Accuracy of Teachers’ Survey Self-Reports

The eight teachers we studied more closely through interview analysis are HQ, NC, MD and EB in Region Z and WH, LS, XN and KN. Those teachers are highlighted in Table 1 and identified in Figures 3 and 4. Region Z teachers generally reported participating in much less curriculum-specific professional development compared to Greene teachers, and they made reference to multiple instructional programs that overshadowed focus on *Everyday Mathematics* in their district. At the same time, the only two Region Z teachers with consistently accurate survey responses – HQ and NC – reported more sustained interaction about their mathematics instruction with other teachers and coaches compared to MD and EB, who over-estimated their practice at least for one year.

By contrast, most Greene teachers experienced intensive curriculum-specific professional development, especially in the first year of their implementation of *Investigations*. That curriculum-specific professional development likely drove more accurate survey responses in Greene. That said, another instructional program overshadowed a focus on *Investigations*-related interaction and learning opportunities in Year 2 of our study in Greene (see Kaufman and Stein 2010 for documentation on that shift). Despite this program shift, XN continued to develop her standards-based practice and moved from “medium” to “high” observation ratings from Year 1 to Year 2, becoming an under-estimator of her practice. Those who did not continue to develop their standards-based mathematics practices in Year 2 – WH, LH, and KN – either sustained accurate reports of “medium” instruction or experienced some drop in their observation ratings or survey reports.

Conclusions: Our analysis leads us to propose the following hypotheses for future controlled studies. First, our study indicates that those teachers with “medium” standards-based practices are more likely to accurately estimate their practices, whereas teachers with “low” practices more often over-estimate their practice and teachers with “high” practices more often under-estimate them. Additionally, our evidence supports the hypothesis that more curriculum-specific interaction and teacher learning opportunities – and an emphasis on fewer other instructional programs – likely lead to more accurate survey reports. Finally, our work indicates that change in teachers’ standards-based practice likely occurs incrementally. These specific hypotheses should be considered in future research, both as possible factors to control for or design against in conducting studies that connect interventions at the teacher level with student outcomes, and as objects of study in their own right, to better understand how to construct and analyze teacher surveys to better reflect actual classroom practice.

Appendices

Appendix A. References

- Coburn, C. and J. L. Russell (2008). District Policy and Teachers' Social Networks. *Educational Evaluation and Policy Analysis*, 30, 203-235.
- Engle, R. A. and R. B. Faux (2006). Fostering Substantive Engagement of Beginning Teachers in Educational Psychology: Comparing two methods of case-based instruction. *Teaching Educational Psychology*, 1, 3-24.
- Hill, H. C. (2005). Content Across Communities: Validating Measures of Elementary Mathematics Instruction. *Educational Policy*, 19, 447-475.
- Hill, H. C., S. G. Schilling, et al. (2004). Developing measures of teachers' mathematics knowledge for teaching. *Elementary School Journal*, 105, 11-30.
- Lampert, M. (1990). When the Problem Is Not the Question and the Solution Is Not the Answer: Mathematical Knowing and Teaching. *American Educational Research Journal* 27, 29-63.
- Kaufman, J.H. and M.K. Stein (2010). Teacher learning opportunities in a shifting policy environment for instruction. *Educational Policy*, 24, 563-601.
- Mayer, D. P. (1999). Measuring Instructional Practice: Can Policymakers Trust Survey Data? *Educational Evaluation and Policy Analysis*, 21, 29-45.
- National Council of Teachers of Mathematics (1989). *Curriculum and Evaluation Standards for School Mathematics*. Reston, VA: Author.
- Ross, J. A., D. McDougall, et al. (2003). A Survey Measuring Elementary Teachers' Implementation of Standards-Based Mathematics Teaching. *Journal for Research in Mathematics Education*, 34, 344-363.
- Shavelson, R. J. and L. Towne, Eds. (2002). *Scientific research in education*. Washington, DC: National Research Council, National Academy Press.
- Shifter, D. (2001). Learning to See the Invisible: What skills and knowledge are needed to engage with students' mathematical ideas? In T. Wood, B. S. Nelson and J. Warfield (Eds.), *Beyond Classical Pedagogy: Teaching elementary school mathematics*. Mahwah, NJ: Erlbaum.
- Spillane, J. P. and J. S. Zeuli (1999). Reform and Teaching: Exploring Patterns of Practice in the Context of National and State Mathematics Reforms. *Educational Evaluation and Policy Analysis*, 21, 1-27.

- Stecher, B., V.-N. Le, et al. (2006). Using Structured Classroom Vignettes to Measure Instructional Practices in Mathematics. *Educational Evaluation and Policy Analysis*, 28, 101-130.
- Stein, M. K. and J. H. Kaufman (2010). Selecting and supporting the use of mathematics curricula at scale. *American Educational Research Journal*, 47, 663-693.
- Stein, M.K. and C. Coburn (2008). Architectures for learning: A comparative analysis of two urban school districts. *American Journal of Education*, 114, 583-626.
- Yin, R. K. (1994). *Case Study Research: Design and Methods, Second Edition*. Thousand Oaks, CA: Sage Publications.

Appendix B. Tables and Figures

Table 1. Participating Teachers, Number of Lesson Observations, and Accuracy of Survey Report Each Year[§]

NA – Not applicable because teacher was not in district that year or did not have sufficient data to be included

District	School	Teacher Initials	Year 1 Lessons	Year 2 Lessons	Year 1	Year 2
Region Z	A	BE	6	6	Accurate	Under estimator
Region Z	A	BH	6	6	Over estimator	Over estimator
Region Z	A	BT	NA	6		Over estimator
Region Z	A	PQ	6	4	Over estimator	Over estimator
Region Z	A	QJ	6	NA	Over estimator	
Region Z	B	DS	6	NA	Over estimator	
Region Z	B	HQ	6	6	Accurate	Accurate
Region Z	B	NC	6	6	Accurate	Accurate
Region Z	B	NH	5	NA	Over estimator	
Region Z	B	OG	NA	5		Accurate
Region Z	B	UF	6	6	Over estimator	Over estimator
Region Z	C	DD	6	NA	Accurate	
Region Z	C	EB	6	6	Over estimator	Accurate
Region Z	C	MD	6	5	Over estimator	Over estimator
Region Z	C	TF	6	NA	Accurate	
Region Z	C	TP	5	NA	Accurate	
Region Z	C	TT	NA	5		Accurate
Region Z	C	UW	6	6	Over estimator	Accurate
Region Z	D	EN	6	6	Accurate	Over estimator
Region Z	D	KD	6	NA	Over estimator	
Region Z	D	KT	6	6	Accurate	Under estimator
Region Z	D	MF	6	NA	Under estimator	
Region Z	D	SD	6	6	Over estimator	Over estimator
Greene	E	BX	6	6	Accurate	Accurate
Greene	E	KE	6	NA	Accurate	
Greene	E	LX	6	NA	Accurate	
Greene	E	QL	NA	6		Accurate
Greene	E	QS	6	6	Under estimator	Accurate
Greene	E	SD	6	NA	Accurate	
Greene	E	SN	6	6	Accurate	Accurate
Greene	F	CD	6	NA	Under estimator	
Greene	F	DT	6	6	Accurate	Accurate
Greene	F	KN	6	6	Over estimator	Accurate
Greene	F	LH	6	6	Under estimator	Accurate
Greene	F	NQ	6	NA	Under estimator	

[§] Teacher initials are pseudonyms. All the teachers above also completed a survey for each year of their participation.

Greene	F	XN	6	6	Accurate	Under estimator
Greene	G	LS	6	6	Accurate	Over estimator
Greene	G	LI	6	6	Over estimator	Accurate
Greene	G	NN	5	6	Over estimator	Accurate
Greene	G	TS	6	6	Accurate	Accurate
Greene	G	WH	6	6	Accurate	Accurate
Greene	H	DN	6	NA	Accurate	
Greene	H	KH	6	NA	Over estimator	
Greene	H	NR	NA	6		Accurate
Greene	H	QK	NA	4		Accurate
Greene	H	SH	6	6	Under estimator	Under estimator
Greene	H	UN	6	6	Accurate	Accurate

Figure 1.

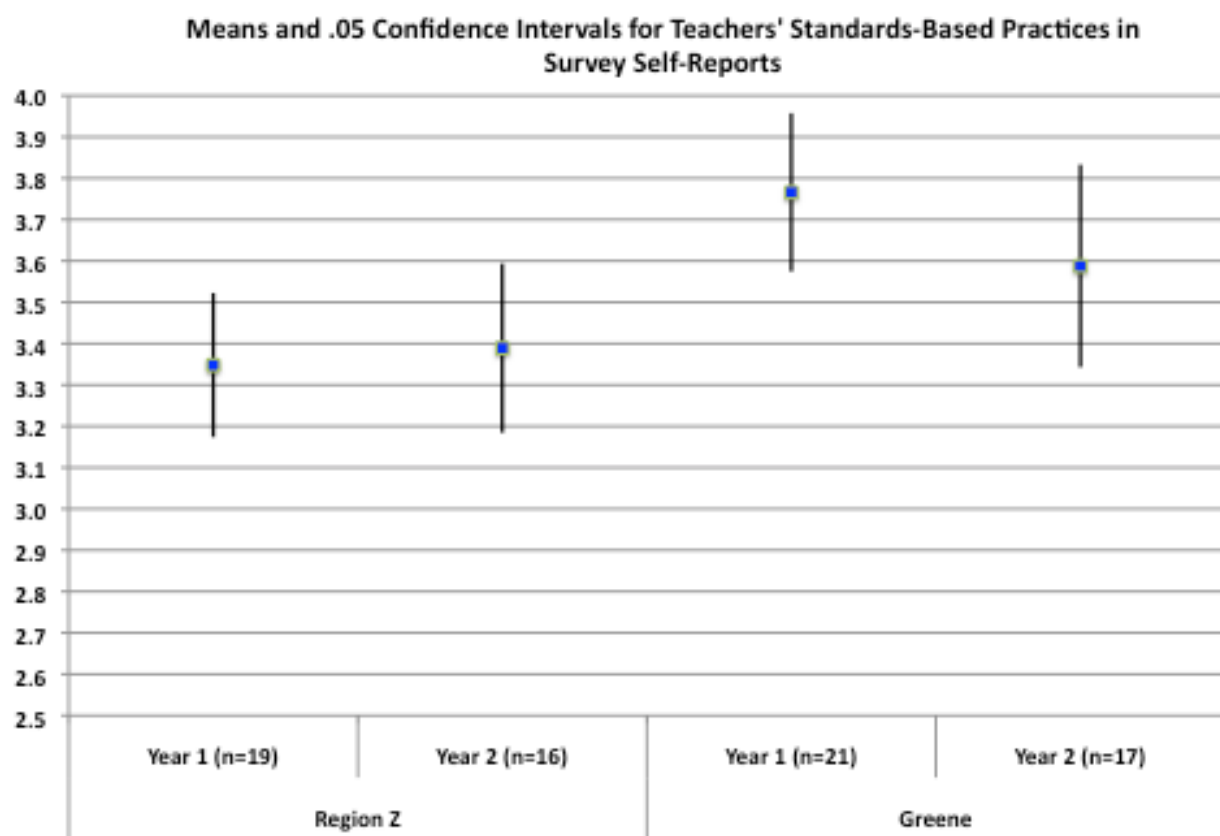


Figure 2.

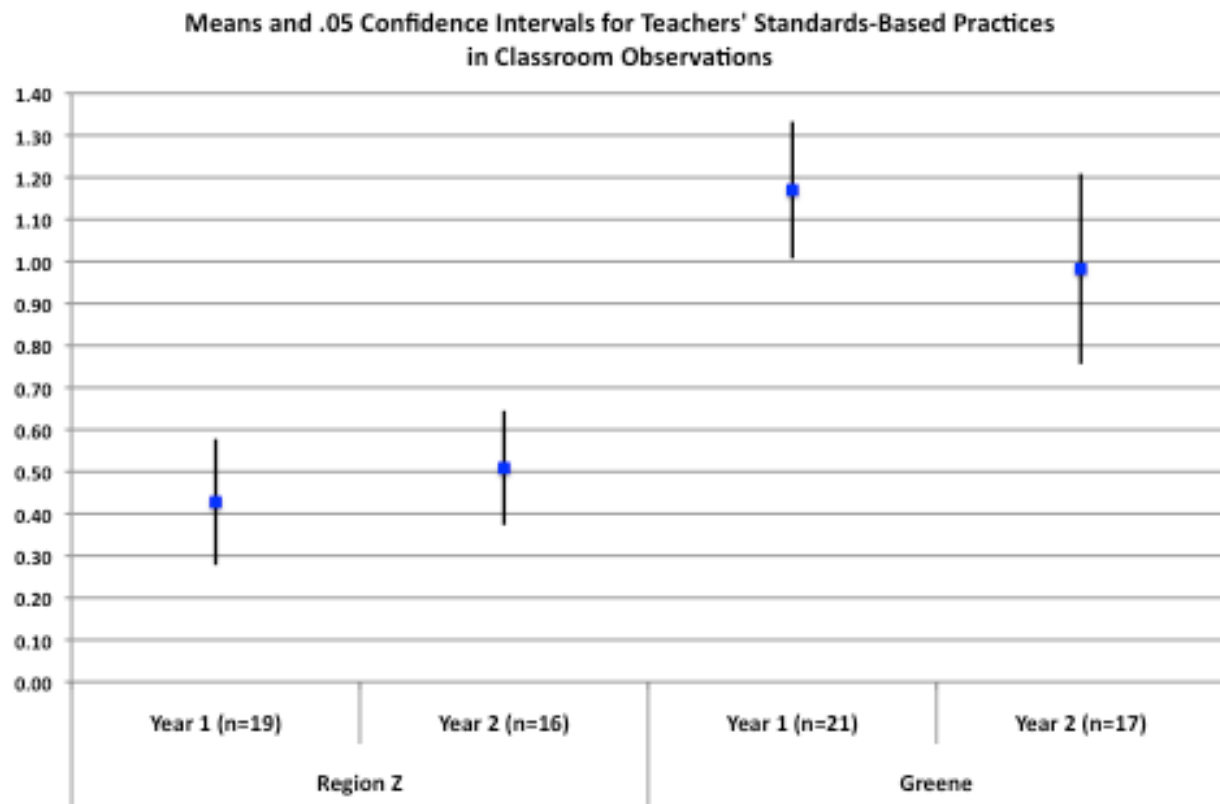


Figure 3.

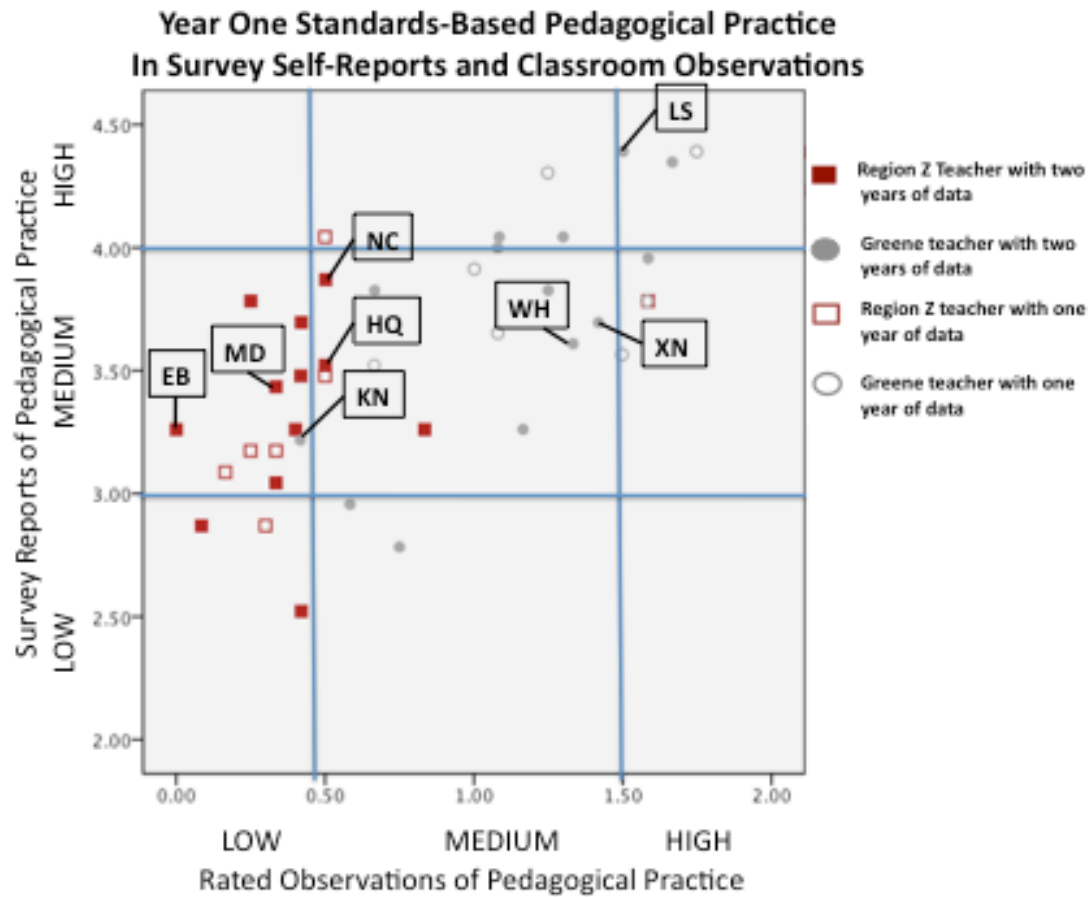


Figure 4.

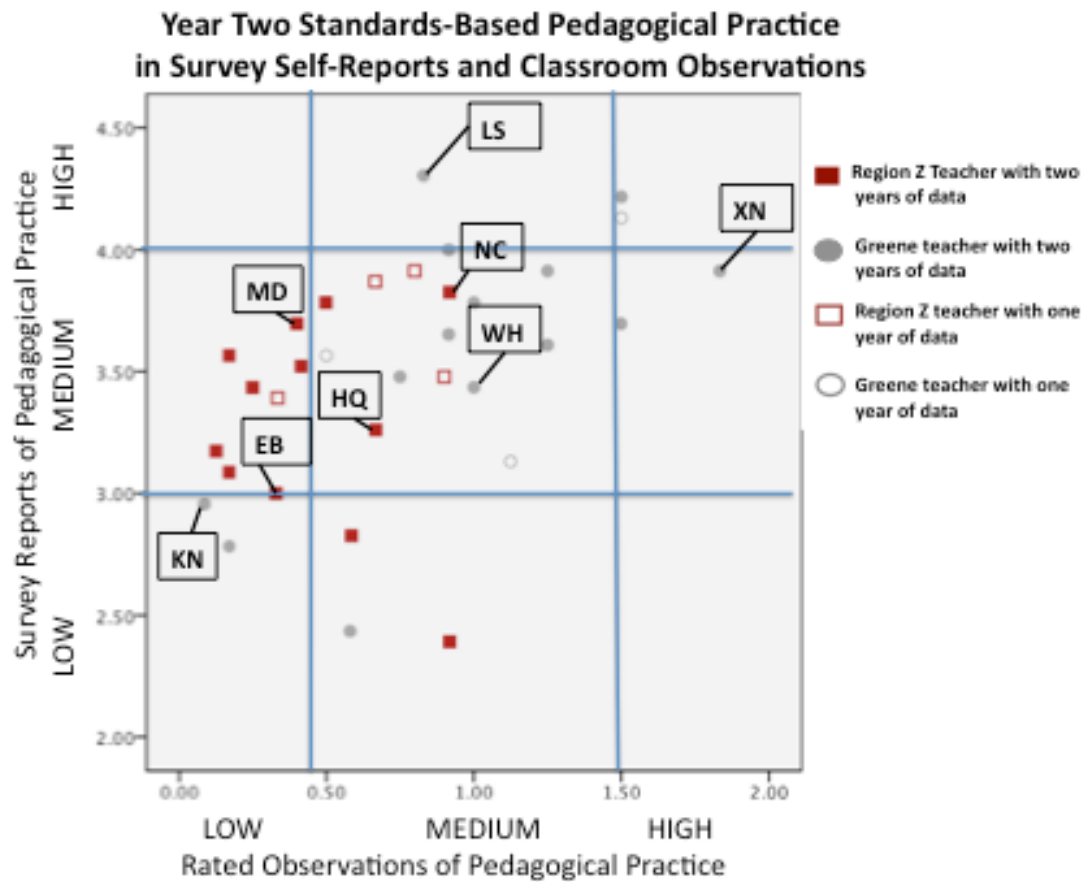


Table 2. Survey and Observation Rating Scales to Measure Standards-Based Practice

Survey items to measure standards-based pedagogical practices (RC=Reverse coded)	
I. The following statements relate to how you teach mathematics. Please answer to what extent you agree or disagree with each statement (scale of 1-5 with 1=Strongly Disagree; 3=Neither Agree nor Disagree; and 5=Strongly Agree)	
1. I like to use math problems that can be solved in many different ways.	
2. When two students solve the same math problem correctly using two different strategies I have them share the steps they went through with each other	
3. I often learn from my students during math time, because my students come up with ingenious ways of solving problems that I have never thought of	
4. When students are working on math problems, I put more emphasis on getting the correct answer than on the process followed (RC)	
5. I don't necessarily answer students' math questions but rather let them puzzle things out for themselves	
6. I like my students to master basic mathematical operations before they tackle complex problems (RC)	
II. How often do you do each of the following in your mathematics instruction? (scale of 1-5 with 1=never, 2=rarely, 3=sometimes, 4=often, 5=always)	
7. pose open-ended questions	
8. engage whole class in discussion	
9. require students to explain their reasoning when giving an answer	
10. ask students to explain concepts to one another	
11. ask students to consider alternative methods for solutions	
III. How you set up work in your classroom (scale of 1-5 with 1=never, 2=rarely, 3=sometimes, 4=often, 5=always) [All items below reverse coded.]	
12. Before assigning a set of problems to my students, I review all procedures/algorithms that can be used to solve the problem. (RC)	
13. Before students begin work on a task, I tell them that they will be able to check the accuracy of their work by checking with me as soon as they've finished. (RC)	
14. When assigning a set of problems, I tell my students which procedure they should use. (RC)	
15. Before turning an open-ended project over to my students, I walk them through an example of how to successfully attack the problem. (RC)	
16. Before turning an open-ended project over to my students, I give them a detailed roadmap to follow through the project. (RC)	
17. I provide students with more steps to follow than what appears in the curriculum that I use. (RC)	
IV. Responding to students (scale of 1-5 with 1=never, 2=rarely, 3=sometimes, 4=often, 5=always)	
18. When students get stuck on a multistep problem, I walk them through the steps they need to perform. (RC)	
19. After students have worked on a particularly challenging assignment, I provide opportunities for them to see how others have approached the assignment.	
20. When students are uncertain about how to get started on an open-ended project, I tell them how to do the first step. (RC)	
21. When a student is unable to complete a task on his/her own, I give him/her a set of steps to follow. (RC)	
22. When students construct their own ways of doing a problem, I have students themselves share their approaches with the rest of the class using their own ways of expressing themselves.	
23. I use students' responses to problems as the fodder for class discussion.	
Prompts for classroom observers to record details about teachers' standards-based pedagogical practices	
Prompts for details about teachers' work to uncover student thinking:	
1. What, if anything, did the teacher do to uncover student thinking? Describe the manner in which the	

<p>teacher provided opportunities for students to make their thinking public. PROVIDE EXAMPLES.</p> <p>2. How did the teacher listen to student thinking? What evidence was there that the teacher tried to understand student thinking? PROVIDE EXAMPLES.</p> <p>3. Describe how the teacher assisted student thinking. To what extent did the teacher help students to identify and articulate the key ideas in their work or thinking? How did she help students represent their thinking and keep track of their work? How did she ask questions that pushed students' thinking? PROVIDE EXAMPLES</p> <p>4. Describe how the teacher made student thinking available for the entire class. Did shared student work come primarily from volunteers or did the teacher appear to have a purpose for whose methods were displayed and in what order? Once students displayed their work, what did the teacher do with it? How did the teacher help students to explain their thinking to the entire class? How did she facilitate class discussions about student work? PROVIDE EXAMPLES</p> <p>5. How, and under what conditions, did the teacher encourage links between students' informal reasoning and more formal, canonical or sophisticated mathematical thinking? PROVIDE EXAMPLES.</p>
<p>Prompts for details about teachers' work to vest intellectual authority in mathematical reasoning:</p> <p>1. What did the teacher expect or allow students to discover on their own and under what circumstances? How to approach problems? The concepts that underlie problems? How to organize and record their work? How to justify their conjectures? PROVIDE EXAMPLES</p> <p>2. What knowledge did the teacher impart or teach to the students and under what circumstances? The steps required to do the mathematical problems? The concepts that underlie the problems? How to attach mathematical notation to their work? The justification of a particular mathematical move? PROVIDE EXAMPLES</p> <p>3. How and by whom was the correctness of a mathematical answer or approach determined? By the answer in the resource materials? By the flawless execution of a procedure? By a calculator? By mathematical logic? PROVIDE EXAMPLES.</p>
<p>Prompts for coders to rate classroom observations (through reading observer's notes and responses to prompts) in regard to 1.) student thinking work and 2.) intellectual authority work</p>
<p>Rate the teachers' work to uncover student thinking based on your reading of the lesson and observers responses to prompts:</p> <p>0 The teacher did no work to uncover student thinking; he/she did most of the talking in the lesson and/or asked questions with short or one-word answers.</p> <p>1 The teacher did some work to uncover student thinking by asking some open-ended questions; by asking for some explanations; by arranging for public sharing of student responses; and/or by listening respectfully.</p> <p>2 In addition to #1 above, the teacher purposefully selected certain students to share their work during whole-class discussion because she wanted the whole class to hear about mathematical approach the student took. However, the teacher <u>did not</u> sequence or connect students' responses in a mathematically meaningful way (i.e. to move the class toward the mathematical goal of the lesson).</p>
<p>Rate the teachers' work to vest intellectual authority in mathematical reasoning based on your reading of the lesson and observers responses to prompts:</p> <p>0 The teacher fostered little or no student construction of mathematical ideas, thinking and/or reasoning. Judgments about correctness were derived from the text or the teacher, with no appeal to mathematical reasoning.</p> <p>1 The teacher fostered some student construction of mathematical ideas, thinking and/or reasoning. However, judgments about correctness were mostly derived from the text or the teacher. Nevertheless, some appeals to mathematical reasoning were made.</p> <p>2 The teacher fostered student construction of mathematical ideas, thinking and/or reasoning. Additionally, judgments about correctness were primarily (most of the time) derived from mathematical reasoning and discussion during the class.</p>